

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-102372

(43) 公開日 平成11年(1999) 4月13日

(51) Int.Cl.<sup>8</sup>

識別記号

F I

G 0 6 F 17/30  
17/27G 0 6 F 15/401  
15/20  
15/403 2 0 A  
5 5 0 F  
3 7 0 A

審査請求 未請求 請求項の数 8 O L (全 9 頁)

(21) 出願番号 特願平9-263318

(22) 出願日 平成9年(1997) 9月29日

(71) 出願人 000005049

シャープ株式会社

大阪府大阪市阿倍野区長池町22番22号

(72) 発明者 池内 洋

大阪府大阪市阿倍野区長池町22番22号 シ  
ャープ株式会社内

(72) 発明者 芥子 育雄

大阪府大阪市阿倍野区長池町22番22号 シ  
ャープ株式会社内

(72) 発明者 黒武者 健一

大阪府大阪市阿倍野区長池町22番22号 シ  
ャープ株式会社内

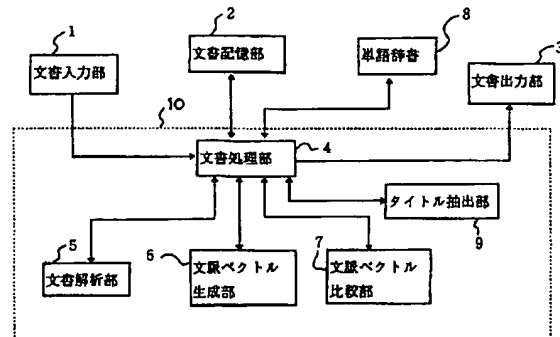
(74) 代理人 弁理士 梅田 勝

(54) 【発明の名称】 文書要約装置及びコンピュータ読み取り可能な記録媒体

(57) 【要約】

【課題】 新聞記事から要約を抽出するには、文書の意味を扱う必要があるが、従来は特定の文書形式や文脈を仮定する必要があった。

【解決手段】 タイトル抽出部9は文書入力部1から入力された文書からタイトルと本文を抽出し、文書解析部5はタイトルを単語に、本文を文および単語に分解する。文脈ベクトル生成部6は、単語辞書8を用いてタイトルおよび本文中の文の文脈ベクトルを生成する。文脈ベクトル比較部7は、タイトルと本文中の各文の文脈ベクトルを比較して各文脈ベクトル間距離を算出する。文書処理部4は各文脈ベクトル間距離を参照して、タイトルに近い複数文を生成して文書の重要部分とする。このように、入力文書を文脈ベクトルを用いて解析することによって、特定の文書形式や文脈を仮定することなく、質の良い重要部分を簡単な処理で抽出できる。



## 【特許請求の範囲】

【請求項1】 入力された文書のタイトルと本文から要約を作成する文書要約装置であって、  
単語の文脈ベクトルが格納された単語辞書と、  
上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、  
上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトルを生成する文脈ベクトル生成手段と、  
上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの距離を算出する距離算出手段と、  
上記算出された距離に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段と、を備えることを特徴とする文書要約装置。

【請求項2】 入力された文書のタイトルと本文から要約を作成する文書要約装置として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体であって、  
単語の文脈ベクトルが格納された単語辞書と、  
上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、  
上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトルを生成する文脈ベクトル生成手段と、  
上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの距離を算出する距離算出手段と、  
上記算出された距離に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項3】 入力された文書のタイトルと本文から要約を作成する文書要約装置であって、  
単語の文脈ベクトルが格納された単語辞書と、  
上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、  
上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトル及び上記文書全体の文脈ベクトルを生成する文脈ベクトル生成手段と、  
上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの第1の距離を算出すると共に、上記文書全体の文脈ベクトルと上記分割された文の文脈ベクトルとの第2の距離を算出する距離算出手段と、  
上記算出された第1の距離に基づいて上記本文から少なくとも一つの文を要約として選出すると共に、上記算出された第2の距離に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段と、を備えることを特徴とする文書要約装置。

【請求項4】 入力された文書のタイトルと本文から要約を作成する文書要約装置として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体

であって、  
単語の文脈ベクトルが格納された単語辞書と、  
上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、  
上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトル及び上記文書全体の文脈ベクトルを生成する文脈ベクトル生成手段と、  
上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの第1の距離を算出すると共に、上記文書全体の文脈ベクトルと上記分割された文の文脈ベクトルとの第2の距離を算出する距離算出手段と、  
上記算出された第1の距離に基づいて上記本文から少なくとも一つの文を要約として選出すると共に、上記算出された第2の距離に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項5】 入力された文書のタイトルと本文から要約を作成する文書要約装置であって、  
単語の文脈ベクトルが格納された単語辞書と、  
上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、  
上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトル及び上記文書全体の文脈ベクトルを生成する文脈ベクトル生成手段と、  
上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの第1の距離を算出すると共に、上記文書全体の文脈ベクトルと上記分割された文の文脈ベクトルとの第2の距離を算出する距離算出手段と、  
上記算出された第1の距離と第2の距離の和に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段と、を備えることを特徴とする文書要約装置。

【請求項6】 入力された文書のタイトルと本文から要約を作成する文書要約装置として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体であって、  
単語の文脈ベクトルが格納された単語辞書と、  
上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、  
上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトル及び上記文書全体の文脈ベクトルを生成する文脈ベクトル生成手段と、  
上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの第1の距離を算出すると共に、上記文書全体の文脈ベクトルと上記分割された文の文脈ベクトルとの第2の距離を算出する距離算出手段と、  
上記算出された第1の距離と第2の距離の和に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項 7】 上記距離算出手段が、上記タイトルと上記分割された文との同一単語の比率を考慮して第 1 の距離を算出するとともに、上記文書全体と上記分割された文との同一単語の比率を考慮して第 2 の距離を算出することを特徴とする請求項 3 または請求項 5 に記載の文書要約装置。

【請求項 8】 上記単語辞書が、単語の特徴ベクトルと共に単語の重要度を格納し、上記距離算出手段が、上記同一単語の比率を算出する際に、上記単語の重要度を用いることを特徴とする請求項 7 に記載の文書要約装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は、電子新聞や情報検索装置等で用いられ、情報の内容を要約して利用者に提示する文書要約装置に関するものである。

【0002】

【従来の技術】電子新聞や情報検索装置は、利用者に情報を提示する装置である。電子新聞は各新聞記事を項目として利用者に提示し、利用者は提示された項目から、自分にとって必要と思われるものを選び、その記事内容を読むことによって情報を得るものである。情報検索装置は、利用者が与えた検索要求に基づいてデータを検索し、検索された各データを項目として利用者に提示し、利用者は電子新聞と同様に、提示された項目から自分にとって必要と思われるものを選び、そのデータ内容を読むことによって情報を得るものである。

【0003】特に新聞記事は、その内容は一般に複数の文からなり、利用者はそのすべてを読むことで、内容を完全に理解することが出来る。従って利用者にとって、記事の文章量が多くなるほど、内容を理解するのに時間と労力が必要とされる。そこで、従来では、利用者の読むべき文章量を減少させるように文章中から要約を抽出することがなされていた。

【0004】ここで、個々の文書の文章量を減少させることによって利用者の手間を軽減するために、文書から要約を抽出する場合を考える。この場合には、文書の文章量を減少させても元の文書の含まれる重要な内容が損なわれないような手法を用いる必要がある。

【0005】従来から提唱されている文書要約の手法としては、主に次の 2 つの手法がある。

【0006】第 1 の手法は、文章を表層的に解析するものである。この手法は、単語の出現頻度解析から文章の重要箇所を決定して元の文書に含まれている単語の組み合わせや文の抽出によって要約文の生成を行うものや、文の文末表現および用語によって文章中における主張文を抽出するものがある。

【0007】第 2 の手法は、文章を意味的に解析するものである。この手法は、事前に文章の形式や文脈を仮定しておいて、その仮定に沿って文章を解析して要約を抽出するものや、文の係り受けの粗密性を用いることによ

って内容の重要性を定義して要約を抽出するものがある。

【0008】また、上述とは全く異なる第 3 の手法として、特開平 6-215049 号公報に示されるように、文脈ベクトルによって、全体の文章と最も意味の近い段落あるいは各段落に最も意味の近い文を求め、それらを要約として提示するものである。

【0009】

【発明が解決しようとする課題】ところが、第 1 の手法と第 2 の手法には各々以下のような問題がある。すなわち、第 1 の手法は、第 2 の手法に比べて簡単に実施できる反面、意味を扱わないので文書中の不要な箇所を重要な箇所と誤って判断してしまうという問題がある。一方、第 2 の手法は、最初の仮定が当てはまらないようなタイプの文書に対しては全く非力であり、内容の重要性の定義自体が困難である上、第 1 の手法に比べて処理が複雑であるという問題がある。

【0010】第 3 の手法は、特定の文書形式や文脈を仮定することなく、簡単な処理によって文書における重要部分を要約として抽出できるが、文書全体あるいは段落全体の文脈ベクトルが、要約を最も良く表すとは限らない。特に新聞記事においては、記事全体よりもタイトルに要約がより良く表現されていると考えられる。

【0011】本発明の目的は、新聞記事のように、タイトルとその本文からより精度の高い要約を求めることのできる文書要約装置を提供することにある。

【0012】

【課題を解決するための手段】請求項 1 に記載の文書要約装置は、入力された文書のタイトルと本文から要約を作成する文書要約装置であって、単語の文脈ベクトルが格納された単語辞書と、上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトルを生成する文脈ベクトル生成手段と、上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの距離を算出する距離算出手段と、上記算出された距離に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段と、を備えることを特徴とする。

【0013】請求項 2 に記載のコンピュータ読み取り可能な記録媒体は、入力された文書のタイトルと本文から要約を作成する文書要約装置として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体であって、単語の文脈ベクトルが格納された単語辞書と、上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトルを生成する文脈ベクトル生成手段と、上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの距離を算出する距離算出手段と、上記算出さ

れた距離に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段として機能させるためのプログラムを記録している。

【0014】請求項3に記載の文書要約装置は、入力された文書のタイトルと本文から要約を作成する文書要約装置であって、単語の文脈ベクトルが格納された単語辞書と、上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトル及び上記文書全体の文脈ベクトルを生成する文脈ベクトル生成手段と、上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの第1の距離を算出すると共に、上記文書全体の文脈ベクトルと上記分割された文の文脈ベクトルとの第2の距離を算出する距離算出手段と、上記算出された第1の距離に基づいて上記本文から少なくとも一つの文を要約として選出すると共に、上記算出された第2の距離に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段と、を備えることを特徴とする。

【0015】請求項4に記載のコンピュータ読み取り可能な記録媒体は、入力された文書のタイトルと本文から要約を作成する文書要約装置として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体であって、単語の文脈ベクトルが格納された単語辞書と、上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトル及び上記文書全体の文脈ベクトルを生成する文脈ベクトル生成手段と、上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの第1の距離を算出すると共に、上記文書全体の文脈ベクトルと上記分割された文の文脈ベクトルとの第2の距離を算出する距離算出手段と、上記算出された第1の距離に基づいて上記本文から少なくとも一つの文を要約として選出すると共に、上記算出された第2の距離に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段として機能させるためのプログラムを記録している。

【0016】請求項5に記載の文書要約装置は、入力された文書のタイトルと本文から要約を作成する文書要約装置であって、単語の文脈ベクトルが格納された単語辞書と、上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトル及び上記文書全体の文脈ベクトルを生成する文脈ベクトル生成手段と、上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの第1の距離を算出すると共に、上記文書全体の文脈ベクトルと上記分割された文の文脈ベクトルとの第2の距離を算出する距離算出手段と、上記算出された第1の距離と第2の距離の和に基づいて上記本文から少なくとも一つの文を

要約として選出する選択手段と、を備えることを特徴とする。

【0017】請求項6に記載のコンピュータ読み取り可能な記録媒体は、入力された文書のタイトルと本文から要約を作成する文書要約装置として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体であって、単語の文脈ベクトルが格納された単語辞書と、上記タイトルを単語に、上記本文を文および単語にそれぞれ分解する文書解析手段と、上記単語辞書を参照して上記タイトルの文脈ベクトル及び上記分割された文の文脈ベクトル及び上記文書全体の文脈ベクトルを生成する文脈ベクトル生成手段と、上記タイトルの文脈ベクトルと上記分割された文の文脈ベクトルとの第1の距離を算出すると共に、上記文書全体の文脈ベクトルと上記分割された文の文脈ベクトルとの第2の距離を算出する距離算出手段と、上記算出された第1の距離と第2の距離の和に基づいて上記本文から少なくとも一つの文を要約として選出する選択手段として機能させるためのプログラムを記録している。

【0018】請求項7に記載の文書要約装置は、請求項3または請求項5に記載の文書要約装置において、上記距離算出手段が、上記タイトルと上記分割された文との同一単語の比率を考慮して第1の距離を算出するとともに、上記文書全体と上記分割された文との同一単語の比率を考慮して第2の距離を算出することを特徴とする。

【0019】請求項8に記載の文書要約装置は、請求項7に記載の文書要約装置において、上記単語辞書が、単語の特徴ベクトルと共に単語の重要度を格納し、上記距離算出手段が、上記同一単語の比率を算出する際に、上記単語の重要度を用いることを特徴とする。

#### 【0020】

##### 【発明の実施の形態】

(実施の形態1) 以下、本実施の形態を図面を用いて説明する。図1に本実施の形態の文書要約装置におけるブロック図を示す。要約作成の対象となる文書が入力される文書入力部1は、キーボードや光学式文字読み取り装置(OCR)、あるいは着脱式外部記憶装置で構成され、処理部10に接続されている。この接続は、通信回線を介して接続されていてもよい。文書記憶部2は、文書入力部1から入力された文書、及び処理部10で生成された要約が格納される。入力文書や要約を出力する文書出力部3は、CRT(カソード・レイ・チューブ)、液晶表示装置(LCD)、プリンタ、あるいは着脱式外部記憶装置で構成されている。この接続は、処理部10に接続され、通信回線を介して接続されていてもよい。

【0021】処理部10内の文書処理部4は、編集や検索等の一般的な文書処理を実施する他に、以下に述べる文書解析部5、文脈ベクトル生成部6、文脈ベクトル比較部7およびタイトル抽出部9を制御して、入力文書の要約を生成する。

【0022】まず最初に、本発明で用いる特徴ベクトルとしての文脈ベクトルについて簡単に説明する。特徴語として、例えば次のような複数の単語「人間，男，女，機械，知識，活動，経験，政治，芸術，科学，…」を用意して特徴空間を定義する。上記特徴語の個数は任意であるが、少なくとも200語～500語程度は用意しておく方が実用上は望ましい。また、特徴語の種類や分野も任意であり、選択に当たっての厳密さは要求されず、特徴が相互にオーバーラップしていても構わない。さらに、要約抽出の対象となる文書の分野が特定の分野である場合には、その分野に特有の特徴語を予め充実させることによって、より品質の高い要約を抽出できることになる。

【0023】単語辞書8には、文脈ベクトルを生成する際に使用される単語を格納し、単語辞書8に格納された各単語と上記特徴語との関連の有無に応じて当該単語を上記特徴空間に配置する。その際における各単語の特徴空間上の位置がその単語の文脈ベクトルであり、この文脈ベクトルは単語に対応させて単語辞書8に格納されている。

【0024】図2に、各単語の文脈ベクトルが定義された単語辞書8の内容の一例を示す。ここで、各要素の配列順序は上述した特徴語の配列順序と同じである。各単語の文脈ベクトルは、単語辞書8内に格納されている単語と各特徴語との関連をその有無によって表現した数字を要素とするベクトルである。すなわち、図2において、関連がある場合には要素“1”を与え、関連がない場合には要素“0”を与えている。ここでは、“1”と“0”を与えているが、これに限定されるものではなく、関連の強度を段階的に表現した数字を要素として与えてもよい。

【0025】図2に例示された単語の文脈ベクトルは以下のことを表現している。すなわち、「人間」という単語は、各特徴語「人間」，“男”，“女”，…とは関連があり、各特徴語「機械」，“知識”，“活動”，“経験”，“政治”，“芸術”，“科学”，…とは関連がないという特徴を表現している。また、「自動車」という単語は、各特徴語「人間」，“男”，“女”，“知識”，“経験”，“政治”，“芸術”，“科学”，…とは関連がなく、各特徴語「機械」，“活動”，…とは関連があるという特徴を表現しているのである。

【0026】本実施の形態において文脈ベクトルを生成する際に用いる単語は、“名詞”および“サ変名詞（語尾に「する」と付けるとサ行変格活用動詞になる名詞）”だけである。したがって、単語辞書8に登録されている単語も名詞およびサ変名詞である。なお、単語辞書8に登録されている単語は名詞およびサ変名詞にしているが、これに限定されるものではない。

【0027】次に、本実施の形態の要約作成処理について図3と図4を用いて説明する。図3は、文書処理部4

による要約作成処理動作のフローチャートである。図4は、電子化されたタイトルとその本文からなるニュース記事である。

【0028】まず、ステップS1で、文書入力部1から要約抽出の対象となる文書が入力されて図4に示す内容すべてが文書記憶部2に記憶される。ステップS2で、タイトル抽出部9によってタイトルと本文が抽出される。タイトル抽出部9は、文書記憶部2から読み出した文書を解析してタイトルと本文とを区別して抽出する。なお、予めタイトルと本文が区別されて入力される場合には、このステップは不要である。一般に電子化されたニュース記事では、予め記事本文とタイトル、およびそれに付随する記事の作成日時、ジャンルなどの情報が一定のフォーマットで付加されている。その場合には単にタイトルに対応する項目と本文に対応する項目から、タイトルと本文とを区別できる。また、タイトルは一般に文の先頭に位置している、あるいはフォントの大きい文字が使われている等のよく知られている特徴を使ってタイトルを抽出し、全文からタイトルを除いたものを本文とするような方法で区別してもよい。以上の方法で抽出されたタイトルと本文は、文書記憶部2の中では区別されて記憶されている。

【0029】ステップS3で、文書解析部5によって、文書記憶部2から本文が読み出されて文単位に分割される。その際に、例えば句点を文の区切りとする。

【0030】ステップS4で、文書解析部5によって、文書記憶部2からタイトルおよび本文が読み出され、この読み出されたタイトルと本文が形態素解析（文章を文法的に解析して単語に分割し、各単語の品詞や活用形等の情報を抽出する手法）によって単語に分解される。図4の例では「ノンバンク整理手続き」がタイトルであり、それが「ノンバンク」「整理」「手続き」の3単語に分解される。文の区切りを句点にして、本文は「日本債券信用銀行の・・・停止した」「一日にも不良債券処理の・・・倒産する」等の文に分割され、前者はさらに「日本」「債券」「信用」「銀行」「の」等の単語に分割され、後者は「一日」「に」「も」「不良」「債券」「処理」「の」等の単語に分解される。このようにして、文書解析部5によってタイトルを単語に、本文を文および単語に分解する。なお、本実施の形態では、単語を抽出する際に形態素解析を行っているが、単語辞書8にある単語との文字列のマッチングによって単語を抽出する方法を用いてもよい。

【0031】そして、得られた単語のうち名詞およびサ変名詞のみが、タイトルおよび上記ステップS3において分割された各文の単位で、文脈ベクトル生成部6に送出される。

【0032】ステップS5で、文脈ベクトル生成部6によって、タイトルの文脈ベクトルと各文の文脈ベクトルとが次のようにして生成される。上記タイトルを構成す

る単語、及び各文を構成する単語の文脈ベクトルは、単語辞書8を引くことによって得られ、さらに、タイトルを構成する単語の文脈ベクトルが加算され正規化されてタイトルの文脈ベクトルが得られ、同様に各文を構成する単語の文脈ベクトルが加算され正規化されて各文の文脈ベクトルが得られるのである。ここで、文脈ベクトルの正規化とは、文脈ベクトルの長さを一定の値に揃えることである。

【0033】この処理を図4を用いて具体的に説明する。タイトルを構成する単語は「ノンバンク」「整理」「手続き」で、それらは名詞またはサ変名詞である。よって、タイトルの文脈ベクトルは、それら3単語の文脈ベクトルを単語辞書8から取り出して加算し正規化して得られる。具体的な計算例として、例えば特徴語が12個で、3単語の文脈ベクトルが  $(0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0)$ 、 $(0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0)$ 、 $(0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0)$  の時、正規化の際に長さを1.0に揃えようとすると、まずそれらを加算したものは、 $(0, 0, 0, 1, 3, 3, 2, 2, 0, 0, 0, 0)$  で、その長さは  $(0^2 + 0^2 + 0^2 + 1^2 + 3^2 + \dots + 0^2)^{0.5} = 27^{0.5}$  である。よって、それを正規化したタイトルの文脈ベクトルは、 $10/27^{0.5} \times (0, 0, 0, 1, 3, 3, 2, 2, 0, 0, 0, 0) \approx (0.0, 0.0, 0.0, 1.9, 5.8, 5.8, 3.8, 3.8, 0.0, 0.0, 0.0, 0.0)$  と求められる。本文についても全く同様にして各文の文脈ベクトルが求められる。

【0034】ステップS6で、文脈ベクトル比較部7によって、上記ステップS5において得られたタイトルの文脈ベクトルと各文の文脈ベクトルとが比較されて各文脈ベクトル間の距離が算出される。その際に算出される2つの文脈ベクトル間の距離は、正規化された当該両文脈ベクトルの内積を用いる。そして、内積値が大きいほど距離が近いのである。つまり、タイトルとの意味が近く、要約としてふさわしい文である。

【0035】ステップS7で、文書処理部4によって、ステップS6において算出されたタイトルと各文との文脈ベクトル間の距離が参照されて、この距離の近い順に本文中の各文が所定数だけ入力文書の要約文として文書記憶部2に格納され、必要に応じて文書出力部3から出力される。文書出力部3は、この要約文のみを出力したり、要約文の箇所をアンダーラインや反転等によって強調された文書全体を出力したりしてユーザの便宜を図ることができる。

【0036】このステップS6とステップS7の処理を具体的に説明する。上述の結果からタイトルの文脈ベクトルは、

$(0.0, 0.0, 0.0, 1.9, 5.8, 5.8, 3.8, 3.8, 0.0, 0.0, 0.0, 0.0)$  であり、  
ステップS5で求められた、本文中の最初の2文の文脈ベクトルが  
 $(0.0, 0.0, 0.0, 6.4, 2.1, 2.1, 6.4, 2.1, 0.0, 0.0, 0.0, 0.0)$

$0, 2.1, 0.0)$

$(0.0, 4.2, 0.0, 0.0, 6.3, 2.1, 0.0, 0.0, 0.0, 6.3, 0.0, 0.0)$

であるとする。タイトルの文脈ベクトルと本文中の最初の文の文脈ベクトルとの内積は、

$0.0 \times 0.0 + 0.0 \times 0.0 + 0.0 \times 0.0 + 1.9 \times 6.4 + 5.8 \times 2.1 + 5.8 \times 2.1 + 3.8 \times 6.4 + 3.8 \times 2.1 + 0.0 \times 0.0 + 0.0 \times 0.0 + 0.0 \times 2.1 + 0.0 \times 0.0 = 68.82$

タイトルの文脈ベクトルと本文中の2番目の文の文脈ベクトルとの内積は、

$0.0 \times 0.0 + 0.0 \times 4.2 + 0.0 \times 0.0 + 1.9 \times 0.0 + 5.8 \times 6.3 + 5.8 \times 2.1 + 3.8 \times 0.0 + 3.8 \times 0.0 + 0.0 \times 0.0 + 0.0 \times 6.3 + 0.0 \times 0.0 + 0.0 \times 0.0 = 48.72$  となる。

【0037】よって、この場合本文中の最初の文の方が2番目の文よりタイトルに近いと判断する。このような内積計算を本文中の全ての文について行い、内積の大きい順に所定数の文が出力されるのである。

【0038】(実施の形態2) 本実施の形態は、実施の形態1と従来の第3の手法とを組合わせて、双方から得られた要約文を共に出力することで、より要約の精度を向上させるものである。

【0039】本実施の形態では、実施の形態1の要約処理に加え、単語辞書8を参照してタイトルと本文の文書全体に含まれる単語から文書全体の意味を表す文脈ベクトルを実施の形態1と同様にして生成し、文脈ベクトル比較部7を用いて文書全体の文脈ベクトルと各文の文脈ベクトルの内積計算を行って距離を求め、距離の近い文ほど文書全体の意味に近いものとなる。双方の手法により求めた距離の順に、タイトルとの意味が近い要約文と、文書全体との意味の近い要約文を共に出力して入力文書の要約となる。

【0040】(実施の形態3) 本実施の形態は、実施の形態2において双方の手法により求めた特徴ベクトル間の距離から総合距離を求めて、総合距離の順に得られた要約文を出力することで、より要約の精度を向上させるものである。

【0041】本実施の形態では、タイトルの文脈ベクトルと各文の文脈ベクトルとの内積値と文書全体の文脈ベクトルと各文の文脈ベクトルとの内積値とを用いることで実現できる。要約を行う文書の各文に文番号を付け、文番号*i*の文脈ベクトルとタイトルの文脈ベクトルとの内積値をST(*i*)、文番号*i*の文脈ベクトルと文書全体の文脈ベクトルとの内積値をSA(*i*)とすると、実施の形態2の2つの方法で計算される距離を総合した総合距離S(*i*)は、

$$S(i) = f(ST(i), SA(i))$$

となる。ここで、*f*は何らかの関数であり、単純に加算もしくはそれらの加重和を取る。具体例としては、

$$S(i) = \alpha \times ST(i) + SA(i)$$

となる。ここで、 $\alpha$ は定数である。このようにして求め

た距離  $S(i)$  の順に、要約文が出力され、入力文書の要約となる。 $\alpha$  の設定は、文書出力部 3 の出力結果を見ながらユーザが文書入力部 1 から設定できるようにしてもよく、また、出力された要約文の所定順位までの  $ST(i)$  の総和と  $SA(i)$  の総和とが均等になるように文書処理部 4 で自動設定するようにしてもよい。このようにすることで、双方の手法による要約が必ず含まれることになる。なお、 $\alpha = 0$  の時は、文書全体の文脈ベクトルとの内積を取った場合と同じになり、 $\alpha \gg 1$  の時はタイトルの文脈ベクトルとの内積値だけで要約文を選ぶ場合と同じになる。

【0042】(実施の形態 4) 本実施の形態は、実施の形態 1 の要約処理に加え、比較する文に含まれている単語の一致する割合を考慮した従来の第 1 の手法を取り入れて、タイトルと各文の距離や文書全体と各文の距離を計算して要約文を出力することで、より要約の精度を向上させるものである。比較する文に含まれている単語の一致する割合の計算は文脈ベクトル比較部 7 で行う。

【0043】具体的には、図 4 のタイトルに含まれる単語は、「ノンバンク」、「整理」、「手続き」の 3 つであり、図 4 の本文の 1 番目の文に含まれるタイトルに含まれている単語は「ノンバンク」の 1 つ、2 番目の文に含まれるタイトルに含まれている単語は「整理」、「手続き」の 2 つ、となっていて、 $i$  番目の文のタイトルに含まれている単語と同一の単語を含んでいる割合  $STK(i)$  は、

$$STK(1) = 1 / 3 \times 100 = 33$$

$$STK(2) = 2 / 3 \times 100 = 66$$

となる。同様に、 $i$  番目の文が、タイトルに含まれている単語と同一の単語を含んでいる割合  $STK(i)$  を計算する。 $i$  番目の文の文脈ベクトルとタイトルの文脈ベクトルとの内積値を  $STV(i)$  とすると、 $i$  番目の文とタイトルとの距離  $ST(i)$  は、

$$ST(i) = f_t(STK(i), STV(i))$$

となる。ここで  $f_t$  は、何らかの関数であり、単純に加算もしくはそれらの加重和を取る。具体例としては、

$$ST(i) = STK(i) + \beta \times STV(i)$$

となる。ここで、 $\beta$  は定数である。なお、 $\beta = 0$  の時は、タイトルに含まれている単語と同一の単語を含んでいる割合だけで要約文を選ぶことになり、 $\beta \gg 1$  の時はタイトルの文脈ベクトルとの内積値だけで要約文を選ぶことになる。

$$STK(1) = 0.8 / (0.8 + 0.5 + 0.6) \times 100 \\ = 42$$

となり、2 番目の文は「整理」、「手続き」の 2 つの単語を含んでいるので、

$$STK(2) = (0.5 + 0.6) / (0.8 + 0.5 + 0.6) \times 100 \\ = 58$$

となる。このように、重要度を取り入れることで、より精度の高い要約を得ることができる。

【0049】以上の実施の形態においては、入力文書を

【0044】また、 $i$  番目の文が、文書全体に含まれている単語と同一の単語を含んでいる割合  $SAK(i)$  を計算する。 $i$  番目の文の文脈ベクトルと文書全体の文脈ベクトルとの内積値を  $SAV(i)$  とすると、 $i$  番目の文のタイトルとの距離  $SA(i)$  は、

$$SA(i) = f_a(SAK(i), SAV(i))$$

となる。ここで  $f_a$  は、何らかの関数であり、単純加算もしくはそれらの加重和を取る。具体例としては、

$$SA(i) = SAK(i) + \gamma \times SAV(i)$$

10 が考えられる。ここで、 $\gamma$  は定数である。なお、 $\gamma = 0$  の時は、文書全体に含まれている単語と同一の単語を含んでいる割合だけで要約文を選ぶことになり、 $\gamma \gg 1$  の時は文書全体の文脈ベクトルとの内積値だけで要約文を選ぶことになる。

【0045】以上のようにして求めた  $ST(i)$  と  $SA(i)$  とを、実施の形態 2 または実施の形態 3 と同様にして距離の順に、要約文が出力され、入力文書の要約となる。なお、 $\beta$  や  $\gamma$  の設定は、文書出力部 3 の出力結果を見ながらユーザが文書入力部 1 から設定できるようにしてもよく、また、出力された要約文の所定順位までの  $ST(i)$  の総和と  $SA(i)$  の総和とが均等になるように文書処理部 4 で自動設定するようにしてもよい。

【0046】(実施の形態 5) 本実施の形態は、実施の形態 4 において、タイトルに含まれている単語と同一の単語を含んでいる割合と文書全体に含まれている単語と同一の単語を含んでいる割合を計算する時に、単語辞書 8 に予め設定した単語の重要度を読み出して、含まれている単語の割合の計算に取り入れるというものである。

【0047】上記  $STK(i)$ 、 $SAK(i)$  は、

$$STK(i) = \sum p_i / \sum p_t \times 100$$

$$SAK(i) = \sum p_i / \sum p_a \times 100$$

と表すことができる。ここで、 $\sum p_i$  は文番号  $i$  の文に含まれている単語の単語の重要度の和を、 $\sum p_t$  はタイトルに含まれている単語の単語の重要度の和を、 $\sum p_a$  は文書全体に含まれている単語の単語の重要度の和を表している。

【0048】具体例を示してみると、図 4 のタイトルに含まれている単語の重要度を、「ノンバンク」を 0.

8、「整理」を 0.5、「手続き」を 0.6 とすると、本文の 1 番目の文は「ノンバンク」の 1 つの単語を含んでいるので、

50 文脈ベクトルを用いて解析することによって、従来の表層的な解析による上記第 1 の手法に比較して、文書における質の良い重要部分を抽出できる。また、従来の文章

13

を意味的に解析する第2の手法に比較して、事前に特定の文書形式や文脈を仮定する必要がないので、入力文書に対する自由度が大きく種々のタイプの文書に適用可能である。さらに、入力文書の構造解析や文脈の意味理解を行って内容の重要性を定義する必要がないので、より簡単な処理によって要約の抽出を行うことができる。また、上記第3の手法に比較しても、タイトルから要約に近い文脈ベクトルが構成されるため、より精度の高い要約文の抽出が可能である。

【0050】なお、上述した要約処理を実行するためのプログラムをフロッピーディスクやCDROM等のコンピュータ読み取り可能な記録媒体に予め記録させておいて、必要に応じてコンピュータにインストールさせて用いてもよい。

【0051】

【発明の効果】以上より明らかなように、本発明によれば、特定の文書形式や文脈を仮定することなく、簡単な処理にて意味的に重要な要約を高精度に生成できる。つまり、入力文書の中で不要な箇所を重要な箇所と誤ったり、仮定した文書形式や文脈に当て嵌まらない入力文書

14

に対して全く非力であったりすることなく、種々のタイプの文書から適切な要約を抽出できる。

【図面の簡単な説明】

【図1】この発明の文書要約装置におけるブロック図である。

【図2】単語の文脈ベクトルが定義された単語辞書の内容の一例を示す図である。

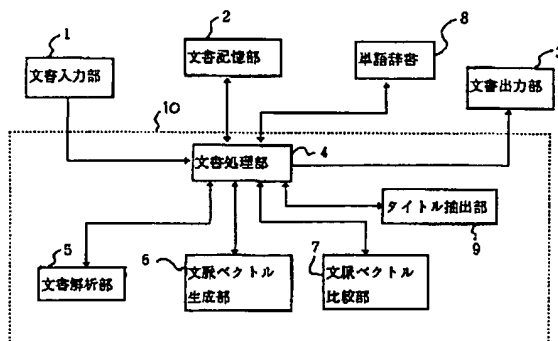
【図3】要約作成処理動作のフローチャートである。

【図4】要約作成処理動作を具体的に説明するための新聞記事の例である。

【符号の説明】

- 1 文書入力部
- 2 文書記憶部
- 3 文書出力部
- 4 文書処理部
- 5 文書解析部
- 6 文脈ベクトル生成部
- 7 文脈ベクトル比較部
- 8 単語辞書
- 9 タイトル抽出部

【図1】



【図2】

単語	人間	自動車	読書	薬品	会社	...	...
人間	1,1,1,0,0,0,0,0,0,...						
自動車	0,0,0,1,0,1,0,0,0,...						
読書	1,0,0,0,1,1,0,0,1,0,...						
薬品	1,0,0,0,1,0,0,0,0,1,...						
会社	1,0,0,0,0,1,1,0,0,0,...						
...	...						
...	...						

【図4】

タイトル： ノンバンク整理手続き

本文： 日本債権信用銀行の関連ノンバンク三社は、銀行などへの元利支払いを停止した。一日にも不良債権処理のため整理手続きに入り、事実上倒産する。整理の方法は特別清算が有力。損失補填は日債銀も他の融資金融機関と同様に融資残高に応じ負担する比例配分方式となる見込み。これを受け、日債銀は海外撤退を.....



【図3】

